

Historiallinen paikkatieto semanttisessa webissä: Biografiasampo.fi

Eero Hyvönen

Semanttinen web tarjoaa lupaavan lähestymistavan tiedon esittämiseen historiallisissa paikannimirekistereissä ja rekisterien julkaisemiseen verkkopalveluina sekä ihmisten että koneiden käyttöön. Esimerkki älykkästä digitaalisesta kulttuurihistoriallisesta sovelluksesta on Suomen kansallisia elämäkertakokoelmia julkaiseva ja rikastava Biografiasampo.fi.

Historiallisen paikkatiedon käsittely on haasteellista: ajan kuluessa paikkojen nimet vaihtuvat, paikkojen sijainti, alue ja tyyppi voivat muuttua ja paikka voi kokonaan hävitä muinaisen Karthagon tapaan. Lisäksi tietomme historiallisista paikosta on epävarmaa (käytiinkö taistelu Vuoksen tuolla tai tällä puolella?), luonteeltaan sumeaa (mistä kohdasta alkaa Korvatunturi?) ja puutteellista (missä on antiikin maantieteilijän Pytheaan tarkoittama Ultima Thule?). Tilannetta ei helpota se, että paikoista on yleensä ollut käytössä monenlaisia ja erikielisiä nimiä.

Historiallisten paikkojen ontologiat ja ontologiapalvelut

Semanttisessa webissä¹ paikannimirekisterit (gazetteer) julkaistaan ns. ontologioina. Niissä jokaiselle paikalle on luotu yksilöivä HTTP URI -tunniste, johon liitetyt ominaisuudet määrittelevät paikan tiedot kuten eri nimet, tyyppin, koordinaatit, laajemman alueen ja paikan osat. Paikat ominaisuuksineen muodostavat verkkomaisen hierarkkisen tietämysgraafin (knowledge graph), joka esitetään semanttisen webin perustana olevan RDF-tietomallin ja standardien avulla².

Olellainen idea on, että HTTP URI -tunniste toimii sekä yksilöivänä tunnisteena että URL-osoitteena, josta käyttäjä voi saada lisätietoa ko. paikasta Tim Berners-Leen linkitetyn datan periaatteiden mukaisesti. Se, mitä tietoa paikkaontologiapalvelusta milloinkin palautetaan, esimerkiksi paikan kuvaava RDF-muotoinen data selaimen JavaScript-sovellukselle tai tästä ihmisen tutkittavaksi muodostettu HTML-sivu, hoidetaan HTTP-protokollan mukaisen sisältöneuvottelun (content negotiation) avulla.

Data rikastuu linkittämällä ja päättelemällä

Ontologiapalveluiden keskeisenä tavoitteena on edistää paikkatietoaineistojen yhteentoimivuutta. Eri data-lähteistä saatava paikkatieto saadaan linkittymään toisiinsa käyttämällä ontologioissa yhteisesti sovittuja tai niihin sillattuja (map, align) tunnisteita.

Esimerkiksi Wikidatan³ ja Wikipedioista louhitun DBpedian⁴ linkitetyn datan palvelut sisältävät valtavan määrän paikkoja, jotka on sillattu muun muassa eri maiden paikannimirekisteristä tuotettuun GeoNames-rekisteriin⁵ ja eri sovellusalojen paikkatietoaineistoihin maailmanlaajuisessa Linked Open Data -pilvessä⁶. Ontologioiden toinen tavoite on datan sisällön rikastaminen datajoukkojen keskinäisellä linkityksellä sekä päätteyllä, joka perustuu semanttisen webin standardien taustalla olevaan formaaliin predikaattilogiikkaan.

¹ E. Hyvönen: Semanttinen web. Linkitetyn avoimen datan käsikirja. Gaudeamus, 2018.

² <https://www.w3.org/standards/semanticweb/>

³ <https://www.wikidata.org/>

⁴ <http://dbpedia.org>

⁵ <https://www.geonames.org/>

⁶ <http://linkeddata.org/>

Kansallisbiografia Biografiasammon ytimessä

Syksyllä 2018 julkistetussa semanttisen webin Biografiasampo.fi-palvelussa⁷ ydinaineistona on Suomalaisen Kirjallisuuden Seuran (SKS) toimittama Kansallisbiografia ja muut 13 100 elämäkertaa, joista on louhittu tekoälyn kieliteknologian ja tietämyksen irrotusmenetelmien (knowledge extraction) avulla 120 miljoonaa yhteyttä sisältävä linkitetyn datan tietämysgraafi⁸.

Paikkaontologian perustaksi otettiin Kansalliskirjaston YSA-sanastosta irrotetuista paikannimistä muodostettu ontologia YSO-paikat⁹, jonka fragmentaarisia hierarkioita muokattiin, täydennettiin ja laajennettiin mm. Sotasampo.fi-palvelun luovutetun Karjalan paikoilla. Paikkaontologia yhdistää monenlaisia Biografiasammossa yhdistyviä, toisiaan rikastuttavia aineistoja. Hierarkkista paikkaontologiaa on hyödynnetty mm. älykkäässä fasettihaussa, jossa vaikkapa ”Saksassa” syntyneitä henkilöitä etsittäessä löytyy Friedenaussa syntynyt sodanaikaisen Päämajan saksalaisen yhdysesikunnan kenraali Waldemar Erfurth, koska Friedenau on osa Berliiniä, joka puolestaan on osa Saksaa.

Järjestelmään sisältyvässä Yhteyshaku-sovelluksessa tavoitteena on päätellä ja löytää henkilöiden ja paikkojen välisiä yhteyksiä ja muodostaa niille suomenkieliset selitykset. Esimerkiksi fasettivalinnoilla ”Italia” ja ”taidemaalari” löytyy tieto, että Elin Danielson-Gambogi vastaanotti Firenzen kaupungin taidepalkinnon vuonna 1899 ja että ”Robert Wilhelm Ekman on luonut vuonna 1844 taideteoksen 'Maisema Subiacosta', joka kuvaa paikkaa Italia”. Jälkimmäisessä tapauksessa yhteys on muodostunut Biografiasampoon yhdistettyjen Ateneumin taidemuseon kokoelmätietojen kautta. Ateneumin lisäksi yhteyshaussa hyödynnetään kokoelmadataa kansallisbibliografia Fennicasta¹⁰, Kirjasammosta¹¹, historiallisia tapahtumia kuvaavasta HISTO-ontologiasta¹² sekä J. V. Snellmanin kootuista teoksista¹³, joista luotiin linkitetyn datan versio.

Visualisointeja paikkaontologialla

Paikkaontologian avulla tietoa voidaan myös visualisoida kartoilla. Mukaan on liitetty historialliset Karjalan kartat ja Senaatin kartasto. Esimerkiksi suomalaisia sotilaita ja pappeja kuvaavat ryhmät kartalla on muodostettu kahdella rinnakkaisella fasettihaulla Biografiasammon elämäkarttojen vertailunäkymässä, jossa elämä kuvataan sinipunaisena nuolena syntymäpaikasta kuolinpaikkaan.

Yhdellä vilkaisulla selviää, että sotilaat liikkuvat pappeja kansainvälisemmin ja vanhemmiten kohti etelää kuten eläkeläiset nykyään. Yhtä kaarta kartalla klikkaamalla pääsee käsiksi kaareen liittyviin elämäkertoihin. Esimerkiksi poikkeava elämänlanka Oulusta Länsi-Siperiaan osoittautuu Siperiaan maanmittaustöiden johtajaksi nimitetyn Gustav Adolf Silverhjelmin aiheuttamaksi.

⁷ <http://biografiasampo.fi/> ja

⁸ <https://seco.cs.aalto.fi/projects/biografiasampo/>

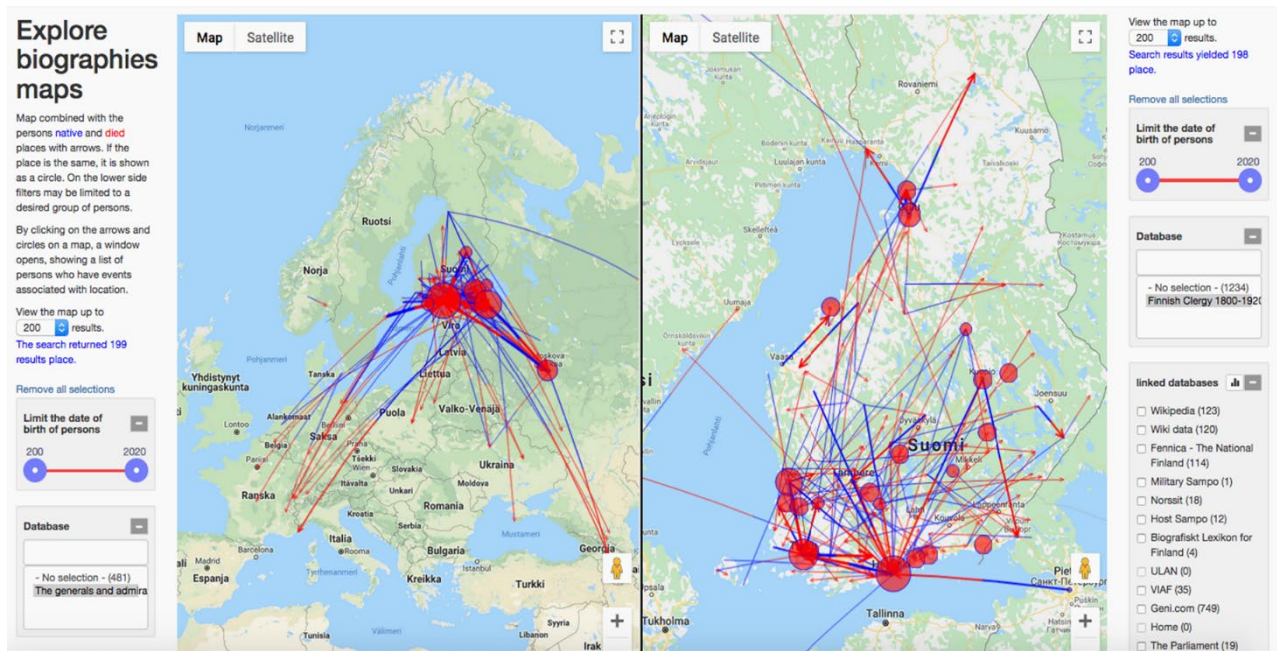
⁹ <https://finto.fi/yso-paikat/fi/>

¹⁰ <https://www.kansalliskirjasto.fi/fi/palvelut/metadatan-muunto-ja-valityspalvelut/avoin-data>

¹¹ <https://www.kirjasampo.fi/>

¹² <https://seco.cs.aalto.fi/ontologies/histo/>

¹³ <http://snellman.kootutteokset.fi/>



Kuva 1. Venäjän sotavoimissa palvelleiden kenraalien ja amiraalien (vasemmalla) ja papiston (oikealla) elämäntien prosopografinen vertailu Biografiasampo.fi-palvelussa.

Haasteena ajallisen muutoksen hallinta

Keskeinen ontologinen haaste historiallisten paikkojen esittämisessä on, että ajan kuluessa paikkojen ominaisuudet muuttuvat. Jos haetaan vaikkapa Suomessa julkaistuja kirjoja, pitäisi Viipurissa painetut kirjat tulla hakutulokseen vain sotia edeltävältä ajalta.

Jo Suomen kuntien lähihistorian ajallinen paikkaontologia on varsin monimutkainen johtuen muutoksista kuntien laajentuessa, yhdistyessä, jakautuessa ja valtorajojen muuttuessa.

Yhden hetkellisen ontologisen kuvauksen, kuten nykyisen Maanmittauslaitoksen Paikannimirekisterin (PNR) sijasta tarvitaan ontologioiden aikasarja, jossa on uusi paikkahierarkia periaatteessa jokaisen muutoksen jälkeen.

Aalto-yliopiston and Geologian tutkimuskeskuksen (GTK) kehittämä Suomen ajallinen paikkaontologia SAPO¹⁴ kuntamuutoksista on yksi tapa lähestyä ongelmaa. Siinä luotiin uusi URI ja identiteetti jokaiselle ajalliselle erilaiselle kunnalle Suomessa 1800-luvulta lähtien aina, kun kunnan alue muuttui. Lisäksi luotiin oma tunniste koko kunnan käsitteelle eri aikoina. Käytettävissä ei ollut juurikaan tietoa kuntien polygoneista eri aikoina vaan ainoastaan kirjoissa dokumentoitua tietoa pinta-aloista ja muutoksista, kuten kuntaliitoksista. Tiedon perusteella muodostetusta ontologiasta voitiin kuitenkin päätellä alueellisia kattavuuksia tiedonhaku varten, kuten että vanhan Viipurin (ennen vuotta 1906) alueesta noin 12 prosenttia kuuluu edelleen Suomeen, nykyiseen Lappeenrantaan.

¹⁴ <https://seco.cs.aalto.fi/ontologies/sapo/>

Suomalaista historiallisten paikkojen ontologiaa ja siihen liittyvää historiallisten karttojen palvelua ollaan tutkimassa ja kehittämässä Aalto-yliopiston ja Helsingin yliopiston Digitaalisten ihmistieteiden keskuksen HELDIG¹⁵ yhteistyönä osana Linked Open Data Infrastructure for Digital Humanities in Finland (LODI4DH)¹⁶ -konaisuutta.

Varteenotettava vaihtoehto ajallisen dimension tuomiseksi paikkaontologiaan on liittää sen ominaisuuksiin lisätietoa siinä, milloin kyseinen ominaisuus on ollut voimassa, esimerkiksi milloin Pietari oli Venäjän pääkaupunki. Näin on menetelty mm. semanttisesti rikkaimmassa kansainvälisessä historiallisten paikkojen rekisterissä, Getty-säätiön Thesaurus of Geographical Names (TGN) -ontologiassa¹⁷, joka nykyisin julkaistaan avoimena linkitettyinä datana verkossa. Suomessa kehitettävissä ontologioissa on hyvä huomioida kansainvälinen yhteentoimivuus paitsi W3C:n yleisiin semanttisen webin standardeihin myös TGN:n kaltaisiin sovellusala-kohtaisiin de facto -standardeihin.

Eero Hyvönen on Helsingin yliopiston Digitaalisten ihmistieteiden keskuksen HELDIG johtaja ja Aalto-yliopiston tietotekniikan laitoksen professori. Hän on julkaissut yli 400 artikkelia ja kirjaa tekoälyyn ja semanttiseen webiin liittyen, uusimpana suomeksi ”Semanttinen web: avoimen linkitetyn datan käsikirja” (Gaudeamus, 2018).

Kotisivu ja yhteystiedot: <https://seco.cs.aalto.fi/u/eahyvone/>

¹⁵ <http://heldig.fi>

¹⁶ <https://seco.cs.aalto.fi/projects/lodi4dh/>

¹⁷ <http://www.getty.edu/research/tools/vocabularies/tgn/>